

# Data Entry Specs for Chinese Text

Wolfgang Schmidle, Martina Siebert,  
Martin Hofmann, Klaus Thoden, Malcolm D. Hyman  
Max Planck Institute for the History of Science, Berlin, Germany

Version 2.1  
16th March 2011

## Contents

<b>1</b>	<b>File Conventions</b>	<b>2</b>
<b>2</b>	<b>General Markup</b>	<b>2</b>
2.1	Page Breaks, Page Numbers and Running Heads . . . . .	2
2.2	Text Blocks . . . . .	3
2.2.1	Headings . . . . .	3
2.2.2	Paragraphs . . . . .	4
2.3	Structured Text . . . . .	6
2.3.1	Tables . . . . .	6
2.3.2	Lists . . . . .	6
2.3.3	Tables of Contents . . . . .	7
2.4	Printed Images . . . . .	8
2.4.1	Figures . . . . .	8
2.4.2	Stamps . . . . .	9
2.5	Unreadable Text . . . . .	9
2.5.1	Characters You are Unsure About . . . . .	9
2.5.2	Unknown Characters . . . . .	10
<b>3</b>	<b>Chinese Characters</b>	<b>11</b>
3.1	General . . . . .	11
3.1.1	Punctuation . . . . .	11
3.1.2	Spaces . . . . .	11
3.2	Character Variants . . . . .	11
3.2.1	Character Variants and Unicode . . . . .	11
3.2.2	Rules for Marking Character Variants . . . . .	12
3.3	Type Styles . . . . .	15
3.3.1	Small Characters . . . . .	15
3.3.2	Underlinings . . . . .	15
3.3.3	Individualized Character Style . . . . .	16
3.3.4	Handwritten Text . . . . .	16
<b>A</b>	<b>Old Radicals That Need to be Marked</b>	<b>17</b>
<b>B</b>	<b>List of All Tags</b>	<b>18</b>

## 1 File Conventions

Save the text in plain text format (.txt) with Unicode utf-8 encoding. If the text is saved in more than one file, number the parts, for example Euclid\_part\_001.txt, Euclid\_part\_002.txt, and so on. Create a zip archive of all files.

Make use of the complete character repertoire found in Unicode version 5.1.

This includes characters in the following Unicode blocks when applicable:

- CJK Unified Ideographs Extension A (U+3400 – U+4DFF)
- CJK Unified Ideographs Extension B (U+20000 – U+2A6DF)

At this point, do not make use of the characters in Extension C or D.

We will also need the list of unknown characters (see section 2.5.2) and the list of character variants (see section 3.2). Please send each list in two versions, namely in the original file format (e.g. RTF, DOC, XLS) and as PDF. If the lists are handwritten, scan them and save them as PDF files.

## 2 General Markup

Type the entire contents of one page, then go on to the next page. Do not mix the contents of different pages.

### 2.1 Page Breaks, Page Numbers and Running Heads

Page breaks are marked by <pb>. If the page has a page number, type it within the <pb> tag, e.g. <pb 六>. Type the page number exactly as it appears in the book. If there is a running head on the page, it is marked by <rh> and </rh>. Type the running head immediately after the <pb> tag.

Type the <pb> and <rh> tags before you type any page content.

The centre section of a traditional printing page (*banxin* 版心) is equivalent to a running head in a western book layout. In this case, repeat <pb> and the running head for each half-page, but add a and b to the page number, e.g. <pb 三a> and <pb 三b>, or <pb a> and <pb b> if there is no page number.

If the characters of the running head are cut off on the scanned page, type them anyway. Type large spaces in the running head as a single IDEOGRAPHIC SPACE character U+3000.

**Please note:** In the digitization of the book, the two half-pages may be on the same scan or on two consecutive scans.

## Examples

<pb 三十二a><rh>泰西事物<起v>原 第十一章</rh>  
<pb 十二a><rh>事物攷辨卷之六十三 <sm>植物</sm> 帶經堂</rh>  
<pb 一a><rh>閩產<錄v>異 卷一<sm>穀屬</sm></rh>  
<pb 一b><rh>閩產錄異 卷一<sm>穀屬</sm></rh>



—> For <sm> see section 3.3.1. An example of two complete half-pages with their running heads can be seen in section 2.2.2.

**Please note:** For < v> (i.e. marking character variants which are not included in Unicode 5.1) see section 3.2.

## 2.2 Text Blocks

Type a return after each line of the printed page.

Do not insert a space at the end of the line.

### 2.2.1 Headings

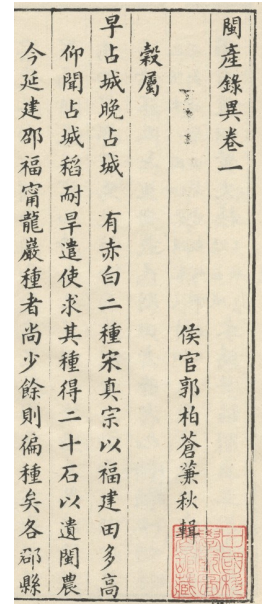
Headings are marked by <h> </h>. If you can identify a heading as the book title, use <ti> for the whole line. If you can identify a heading as the name of the author, compiler, proofreader etc. (*ti* 題), use <ti> too. If a text has different levels of headings, use <h 1> for all headings on the highest level, <h 2> for all headings on the second highest level, and so on.

If a heading, book title, author, etc. is indented, do not mark this. Type large spaces in the heading as a single IDEOGRAPHIC SPACE character U+3000.

## Example

```
<ti>閩產<錄v>異卷一</ti>
<stamp>
<ti>侯官郭柏<蒼R><兼R>秋輯</ti>
<h 1>穀屬</h>
<h 2>早占城晚占城</h> <p>有赤白二種宋真宗以福建田多高
仰聞占城稻耐旱遣使求其種得二十石以遺閩農
今延建邵福甯龍巖種者尚少餘則徧種矣各郡<縣v>
... </p>
<h 2> ... </h> <p> ...
... </p>
```

→ For <stamp> see section 2.4.2. For <p> see section 2.2.2.



**Please note:** In the transcription of this example, there are three headings on two different levels (<h 1>, <h 2>, <h 2>). In the example in section 2.2.2 there are two headings, both on level 2 (<h 2>, <h 2>). <h 1> appears on other pages of the book.

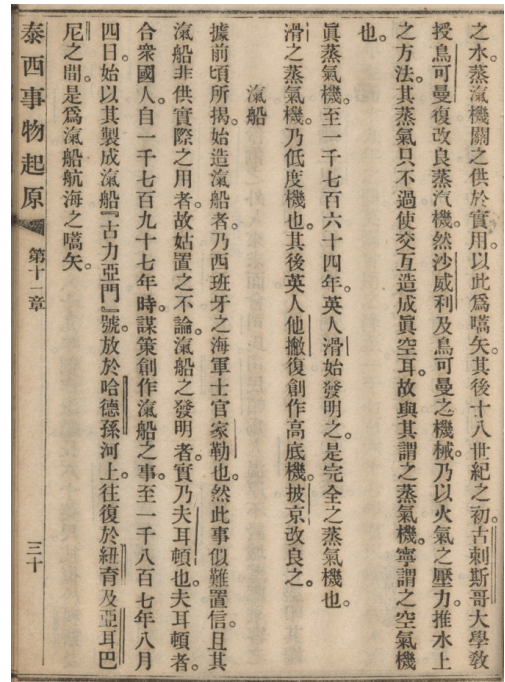
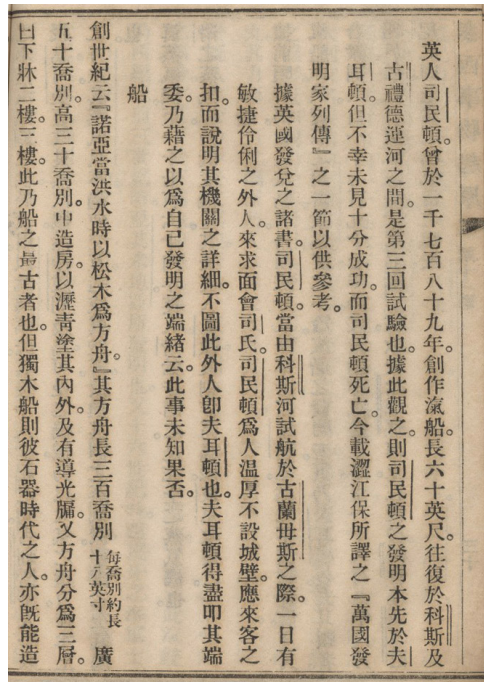
## 2.2.2 Paragraphs

Paragraphs are marked by <p> and </p>. Indented paragraphs are marked by one or more *i*, e.g. <p *iii*> for a paragraph that is indented by three character spaces. Out-dented paragraphs are marked by one or more *x*, e.g. <p *x*>.

Every part of the book has a base line where all unindented paragraphs start. The base lines in the preface or table of contents may be different from the base line in the main text. Mark indentations relative to the base line. The indentation symbols *i* and *x* always refer to the first line of the paragraph. The remaining lines of the paragraph may have the same or a different indentation, which is not marked. If a paragraph is preceded by a sub-heading in the same line, as in the example in section 2.2.1, do not mark the indentation at all.

Make sure that for each <p> there is a corresponding </p> somewhere. If a paragraph starts and ends on different pages, the <p> and </p> tags are on these different pages.

## Example



<pb 三十a><rh>泰西事物<起V>原 <sm>第十一章</sm></rh>

之水。蒸氣機關之供於實用。以此爲嚆矢。其後十八世紀之初。古刺斯哥大學教授鳥可曼復改良蒸汽機。然沙威利及鳥可曼之機械。乃以火氣之壓力。推水上之方法。其蒸氣只不過使交互造成真空耳。故與其謂之蒸氣機。寧謂之空氣機也。</p>

<p>真蒸氣機。至一千七百六十四年。英人滑始發明之。是完全之蒸氣機也。滑之蒸氣機。乃低度機也。其後英人他撒復創作高底機。披京改良之。</p>

<h 2>汽船</h>

<p>據前頃所揭。始造汽船者。乃西班牙之海軍士官家勒也。然此事似難置信。且其汽船非供實際之用者。故姑置之不論。汽船之發明者。實乃夫耳頓也。夫耳頓者。合衆國人。自一千七百九十七年時。謀策創作汽船之事。至一千八百七年八月四日。始以其製成汽船『古力亞門』號。放於哈德孫河上。往復於紐育及亞耳巴尼之間。是爲汽船航海之嚆矢。</p>

<pb 三十b><rh>泰西事物起原 <sm>第十一章</sm></rh>

<p i>英人司民頓。曾於一千七百八十九年。創作汽船。長六十英尺。往復於古刺斯哥及古禮德運河之間。是第三回試驗也。據此觀之。則司民頓之發明本先於夫耳頓。但不幸未見十分成功。而司民頓死亡。今載澀江保所譯之『萬國發明家列傳』之一節以供參考。</p>

<p iii>據英國發兌之諸書。司民頓。當由科河試航於古蘭毋斯之際。一日有敏捷伶俐之外人。來求面會司氏。司民頓爲人溫厚不設城壁。應來客之扣。而說明其機關之詳細。不圖此外人即夫耳頓也。夫耳頓得盡叩其端委。乃藉之以爲自己發明之端緒云。此事未知果否。</p>

<h 2>船</h>

<p>創世紀云『諾亞當洪水時以松木爲方舟』。其方舟長三百喬別。每喬別約長十八英寸。廣五十喬別。高三十喬別。中造房。以瀝青塗其內外。及有導光牖。又方舟分爲三層。曰下牀二樓。三樓。此乃船之最古者也。但獨木船則彼石器時代之人。亦既能造

—> For <sm> see section 3.3.1. For <sl> and <dl> see section 3.3.2.

**Please note:** For the marking of character variants (< V>, < R>, < RV>) see section 3.2. Only the first appearance of each character variant is marked. For instance, 起 is marked in the first running head but not in the second one. (In addition, every

example in the Data Entry Specs is treated as if it was the beginning of a text, i.e. 起 is marked on the first half-page even though it already appears on earlier pages of the text.)

## 2.3 Structured Text

### 2.3.1 Tables

A table is marked by `<tb>` and `</tb>`. Use # as field separators. Type a return after each row. Do not type horizontal or vertical lines.

Type the `<tb>` and `</tb>` tags on separate lines. Do not mark indentations. The field separators may be lines or large spaces.

If you can identify a single space within a name etc. as a decorative space to make the table layout optically more pleasing, do not type it.

#### Example

`<p>`今日諸國所用文字之數如左`</p>`

`<tb>`

英吉利	#	二十六
法蘭西	#	二十三
西班牙	#	二十七
希臘	#	二十四
斯格拉	<code>&lt;窩v&gt;</code>	尼亞 # 二十七
德意志	#	二十六
意大利	#	二十
俄	<code>&lt;羅v&gt;</code>	斯 # 四十一
拉丁	#	二十三
希伯流	#	二十二
梵字	#	五十

`</tb>`

### 2.3.2 Lists

A list is marked by `<list>` and `</list>`. Use # for large spaces, if there are any.

Type the `<list>` and `</list>` tags on separate lines. If the items on consecutive text lines belong to the same list entry, use # at the beginning of the next line. Do not mark indentations.

Unlike in tables, type each single space, i.e. do not omit single spaces even if they seem to be merely decorative.



## Example

```
<list>
  (第一) 意大利語 # 西班牙語 # 法蘭西語 等
# 同由拉丁語出
  (第二) 俄<羅V>斯語 # 波蘭語 # 波希密亞語 等
# 同由斯拉夫語出
  (第三) 威爾斯語 # 希臘語 # 不列顛語 等
# 同由塞爾語出
</list>
<p>以上屬於亞利安語系</p>
<list>
  (第四) 亞刺比亞語 # 希伯流語 # 叙利亞語 等
# 同由塞美的語出
</list>
```

(第一)	意大利語	西班牙語	法蘭西語 等
(第二)	俄羅斯語	波蘭語	波希密亞語 等
(第三)	威爾斯語	希臘語	不列顛語 等
以上屬於亞利安語系			
(第四)	亞刺比亞語	希伯流語	叙利亞語 等
同由塞美的語出			

→ For an example of a list-like structure without large spaces see section 2.3.3 (second example).

## 2.3.3 Tables of Contents

A table of contents is marked by `<toc>` and `</toc>`. If the table of contents has a table-like structure or a list-like structure with large spaces, use #.

Type the `<toc>` and `</toc>` tags on separate lines. A table of contents may look like a table or like a list.

### Examples

a table-like table of contents

a list-like table of contents  
without large spaces

```
<toc>
<h>第一章 天時</h>
日月 # 日月<蝕R> # 地球 # 地球之圓體
地動說 # 遊星 # 七曜日 # 晝夜
時間 # 年月 # 歲首 # <紀R>元
三時代 # 天氣豫報
<h>第二章 地理</h>
亞美利加 # 奧斯土刺利亞 # 蘇彝士河 # 山
堤埭 # 橋 # 周航地球 # 大洪水
...
```

```
<toc>
...
安南稻 米麥 牛尾粟<sm>雀粟 鷺掌粟 狗尾\\粟 虎尾粟 <黃V>粟</sm>
黃粱 釣鈎黍<sm>馬尾黍 番黍 鴨脚黍\\ 黑黍 長<芒R>黍 膏黍</sm> 豆<sm>白豆\\ 黃
豆 黑豆 綠豆 豇豆 豌豆 赤小豆 青豆\\ 褐豆 刀豆 虎爪豆 蠍眼豆 皂<莢R>豆 羊
... </sm>
</toc>
```

第一章 天時	日月	地動說	時間	三時代	第二章 地理	亞美利加	堤埭
日月蝕	遊星	年月	天氣豫報	奧斯土刺利亞	蘇彝士河	山	大洪水
地球	七曜日	歲首	紀元	周航地球			

安南稻	米麥	牛尾粟	雀粟	鷺掌粟	狗尾	虎尾粟	黃粱	釣鈎黍	馬尾黍	番黍	鴨脚黍	黑黍	長芒黍	膏黍	白豆	黃豆	黑豆	綠豆	豇豆	豌豆	赤小豆	青豆	褐豆	刀豆	虎爪豆	蠍眼豆	皂莢豆	羊豆
-----	----	-----	----	-----	----	-----	----	-----	-----	----	-----	----	-----	----	----	----	----	----	----	----	-----	----	----	----	-----	-----	-----	----

## 2.4 Printed Images

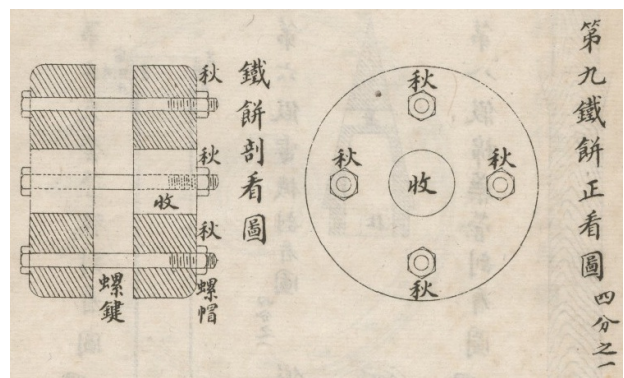
### 2.4.1 Figures

Where a figure occurs in the text, type `<fig>` on a separate line. If you can identify a caption of the figure, mark it by `<cap>` `</cap>`. Additional text that describes parts of the figure is marked by `<desc>` `</desc>`. Use a single `<var>` `</var>` tag for variable names and numbers. Finally, type a closing `</fig>` tag on a separate line.

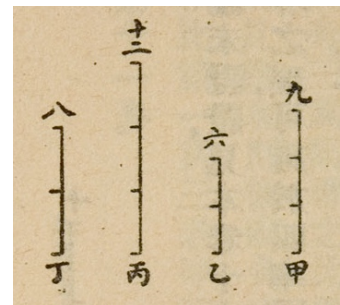
Type all `<cap>`, `<desc>` and `<var>` tags, in this order, on separate lines. A figure may have more than one description. If the same description is repeated in a figure, type it only once. Separate the variable names and numbers in the `<var>` `</var>` by spaces. Type `<fig/>` instead of `<fig>` and `</fig>` to mark simple figures without `<cap>`, `<desc>` or `<var>` tags.

#### Examples

```
<fig>
<cap>第九鐵餅正看圖<sm>四分之一</sm></cap>
<desc>秋</desc>
<desc>收</desc>
</fig>
<fig>
<cap>鐵餅剖看圖</cap>
<desc>秋</desc>
<desc>螺帽</desc>
<desc>收</desc>
<desc>螺鍵</desc>
</fig>
```



```
<fig>
<var>九 甲 六 乙 十二 丙 八 丁</var>
</fig>
```



```
<pb a><rh>泰西事物<起v>原 <sm>第三章</sm></rh>
...
<fig>
<desc>第二大派</desc>
<desc>希伯流語</desc>
<desc>亞刺比亞語</desc>

<pb b><rh>泰西事物起原 <sm>第三章</sm></rh>
<desc>非尼西亞語</desc>
</fig>
...
```





<p>西洋大彈式十種 凡彈必合銃口徑以爲圓形故不預定大小斤數</p>

<fig>

<desc>圓彈</desc>

</fig>

<fig>

<desc>中空<迎R>風  
其聲如雷</desc>

<desc>響彈</desc>

</fig>

<fig>

<desc>中用百鍊鋼條兩頭銼  
尖鑄時先定中<線R>毋使  
稍偏長短致  
有輕重低昂  
不能直貫</desc>

<desc>遇賊攻

寨勢如

拉朽</desc>

</fig>

<fig>

<desc>彈形兩分中<縮R>百鍊  
鋼條不拘長短點放

<進R>發橫拉如火龍</desc>

<desc>鍊彈</desc>

</fig>

<fig>

<desc>最厚之城用十<餘R>彈先  
鑿破磚石<繼R>以員彈推  
倒</desc>

<desc>攻城</desc>

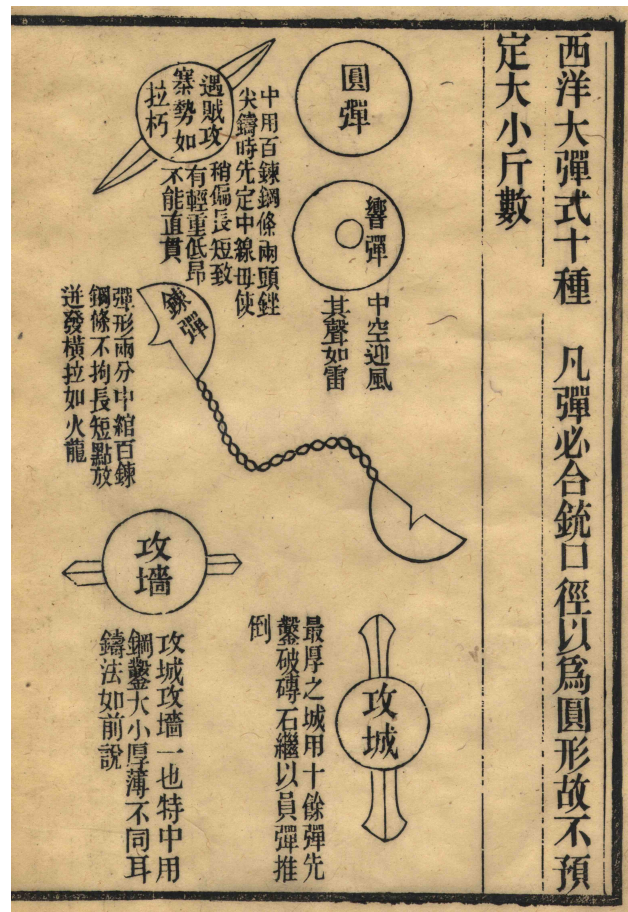
</fig>

<fig>

<desc>攻城攻牆一也特中用  
鋼鑿大小厚薄不同耳  
鑄法如前說</desc>

<desc>攻牆</desc>

</fig>



## 2.4.2 Stamps

Stamps are marked by <stamp>. Type the <stamp> tag on a separate line. Do not type the text in the stamp.

→ For an example see section 2.2.1.

## 2.5 Unreadable Text

### 2.5.1 Characters You are Unsure About

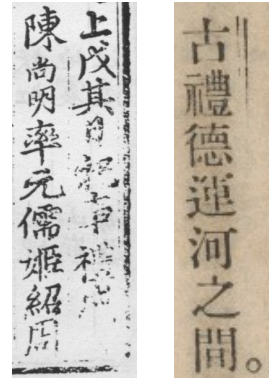
If you are not sure about a character, type <?> after it. If you are unsure about a whole paragraph, type <?> directly after the <p> tag, i.e. <p><?>. A completely unreadable character is typed as @. If many characters are unreadable, use <gap> instead of @.

Use one @ for each unreadable character, e.g. unr@@dable. If in doubt, use <gap>, e.g. unr<gap>dable. If you are unsure about a group of characters, for example a whole word, do not type <?> repeatedly for every character, e.g. type word<?> rather than w<?>o<?>r<?>d<?>.

### Examples

上戌其日祀事禮成<?>@  
陳尚明率元儒姬紹周<?>

<d1>古禮</d1><?>德<運R>河之間。



**Please note:** In the second example, the characters are readable but the double line (see section 3.3.2) is badly printed.

→ For unknown rather than unreadable characters please refer to section 2.5.2.

### 2.5.2 Unknown Characters

If there is an unknown character in the text, i.e. a character variant which is readable but where you cannot identify the standard character, add it to the numbered list of unknown characters. From then on, type its number whenever it occurs in the text, e.g. <001>.

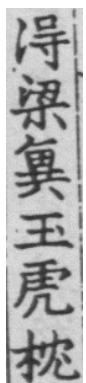
Before you create a number for an unknown character, first check whether it is already on the list of unknown characters. Assign the number <001> to the first unknown character, <002> to the second unknown character, and so on. Do not assign the same number twice. Use this number to type the unknown character. Always use the same number if the same unknown character occurs again.

**Please note:** Make sure that for a given text there is a single list containing all unknown characters, and that everyone uses this list. When the text is sent back to us, we will need a copy of this list. (See also section 1.)

### Example

*an unknown character*

<得v>梁<001>玉<虎v><枕v>



→ For unreadable rather than unknown characters please refer to section 2.5.1. For < v>, i.e. if a character variant is not included in Unicode but you can identify the standard character, see section 3.2.

**Please note:** <001> is actually an unusual version of the variant 冀 of the character 冀. If you can identify the standard character 冀, mark it as <冀v> instead of <001>.

If in addition you can identify it as a version of the character variant 冀 (i.e. as a version of a character variant with a separate Unicode codepoint), mark it as <冀v>. (This special case is not mentioned in section 3.2.)

### 3 Chinese Characters

#### 3.1 General

##### 3.1.1 Punctuation

Type the punctuation to the right of characters.

→ For an example see section 2.2.2.

##### 3.1.2 Spaces

Type spaces in Chinese text as the IDEOGRAPHIC SPACE character U+3000.

In running heads (<rh>, see section 2.1) and headings (<h> and <ti>, see section 2.2.1), type large spaces as a single ideographic space character. In tables and lists (<tb>, <list> and <toc>, see section 2.3), use # for large spaces.

If you encounter a large space in a normal paragraph (<p>, see the example in section 2.2.2, third line from the bottom), make sure that none of the cases above apply. If it is indeed a normal paragraph, type the large space as more than one ideographic space, according to its length.

→ For an example of large spaces that are typed as a single ideographic space character, see section 2.1. For an example of large spaces that are typed as #, see section 2.3.1.

#### 3.2 Character Variants

##### 3.2.1 Character Variants and Unicode

If a character variant is included in Unicode 5.1, type it. Do not normalize the variant.

For example, if the character variant 歷 of the character 歷 occurs in the text, type 歷 (U+6B74). Rather than working with the reference glyphs of the Unicode codepoints, you may use the fonts Sun-ExtA and Sun-ExtB normatively.

If a character variant is not included in Unicode 5.1, type the standard character instead. In addition, mark it by < R> or < V> if the rules in section 3.2.2 apply.

→ If you cannot identify the standard character, treat it as unknown character (see section 2.5.2).

If a character variant should be marked, add it to the Character Variants List, i.e. add the image and how you have marked it. In addition, mark the first occurrence of the character variant in the text. After its first occurrence, silently normalize it. Begin a new Character Variants List for each text.

**Please note:** In rare cases, the same text may contain the standard character and one or more different variants of this character. Please make a note in the Character Variants List, e.g. “S + V” or “V1 + V2” or “S + V1 + V2”. Include an image of each variant. In the text, mark the first occurrence of each variant using < V1>, < V2>, etc. Do not mark the standard character.

### 3.2.2 Rules for Marking Character Variants

This section contains **Rule S** for simple differences that need not be marked, **Rule R** for marking old radicals, and **Rule V** for other variations.

**Rule S (silently normalizing simple differences):** If a character variant differs from the standard character only in the following points, type it as the standard character and do not mark the difference.

In paragraphs with individualized character style (see section 3.3.3) and in handwritten text (see section 3.3.4), silently normalize character variants if they are not in Unicode.

*List: simple differences (with excerpts from ISO 10646)*

#### a) Differences in rotated strokes/dots

- Do NOT mark 𠂇·𠂇, 勹·勹, 羽·羽, 酋·酋
- Do NOT mark 之 (之) where the dot at the top touches or even overshoots the horizontal stroke. Similarly, do not mark the difference between 斥 and 斥.
- + DO mark the difference between 亠 and 十 and between 卜 and 十.  
Examples: 懷 <懷V>, 肺 <肺V>.

#### b) Differences in overshoot at the stroke initiation and/or termination

- Do NOT mark 身·身, 雪·雪
- Do NOT mark 瓦 (瓦), 耶 (耶), or 承 (承). Also do not mark 𠂇 (𠂇), i.e. the differences in the overshoot/non-overshoot of the character part 夕.
- + DO mark vertical (豎) and left-slanted (撇) overshoots.  
Examples: 割 <割V>, 除 <除V>, 楔 <楔V>, 鄂 <鄂V>, 蛇 <蛇V>.

c) Differences in contact of strokes

- Do NOT mark 奧·奧, 酉·酉, 兒·兒, 查·查
- + DO mark connections of the vertical strokes of the component 日 with a horizontal stroke 一 below that look like a change from 旦 to 且.  
Example: The component 易 in 盪 <盪V>.

d) Differences in protrusion at the folded corner of strokes

- Do NOT mark 巨·巨
- Do NOT mark 繩 (繩) or 斲 (斲).
- Do NOT mark the difference between the character components 己 (ji “self”) and 巳 (yi “already”) unless you are completely sure.
- + DO mark the difference between 巳 and 己 / 己.  
Examples: 忌 <忌V>, 包 <包V>.
- + DO mark the difference between 巳 and 巳 and the difference between 皿 and 𠂔.  
Examples: 卷 <卷V>, 服 <服V>, 衆 <衆V>.

e) Differences in bent strokes

- Do NOT mark 西·西
- + DO mark the difference between ㄹ and ㄹ and the difference between 朮 and 朮.  
In general, mark differences in the shape of a stroke if they lead to a different character component.  
Examples: 兪 <兪V>, 述 <述V>.

f) Differences in folding back at the stroke termination

- Do NOT mark 朱·朱

g) Differences in accent at the stroke initiation

- Do NOT mark 父·父, 丈·丈, 乚·乚

h) Differences in “rooftop” modification

- Do NOT mark 八·八, 宀·宀



- Do NOT mark for example 肉 (肉).
- + DO mark the difference between 几 and 儿.  
Examples: 兜 <兜V>, 虎 <虎V>.

i) Straightforward combinations of the above differences

- Do NOT mark 刃・刃・刃

**Please note:** The rules a) to i) leave the number of strokes unchanged. Consequently, changes in the number of strokes cannot be silently normalized and must be marked.

Note, however, that in some printing styles the strokes 丩 (*shuti* 豎提) and 乚 (*shuzhe* 豎折) appear to be a combination of two strokes, i.e. “vertical” (丨) plus “upward-slanted” (丿). Do not mark these printing styles.

Examples: 裏 (裏), 震 (震), 改 (改), 粵 (粵), 仰 (仰).

**Rule R (marking old radicals):** If a character such as 絨 (U+7D68) exists in Unicode 5.1 only with the new shape 糸 (U+7CF9) of the radical (*bushou* 部首) and not with the old shape 糸 (U+7CF8), mark the character variant with the old radical by < R>, for example <絨R>. Apply this rule only if the difference between the old and the new radical is not covered by rule S.

→ Appendix A contains a list of radicals that need to be marked.

**Please note:** Rule R also applies when a character component from the sound-part of the character has swapped places with the actual radical.

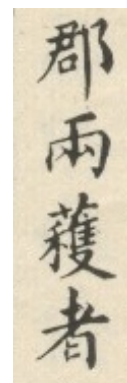
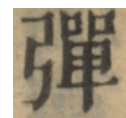
Examples: 穫 <穫R>, 墜 <墜R>.

**Rule V (marking other variations):** If the non-radical part of the character variant differs from that of the standard character, mark the character variant with < V>. Apply this rule only if the difference is not covered by rule S.

*Examples*

<彈V>

郡兩<穫R>者



If both rule R and rule V apply, mark the character variant by < RV>; see the example <過RV> in section 2.2.2.

### 3.3 Type Styles

#### 3.3.1 Small Characters

Strings of small characters are marked by `<sm>` `</sm>`. Indicate half-column breaks by `\\`.

→ For an example see section 2.3.3 (second example). This examples includes strings of small characters over more than one line of text.

#### 3.3.2 Underlinings

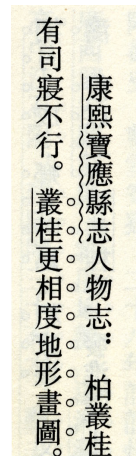
A single line next to characters is marked by `<sl>` `</sl>`. A double line next to characters is marked by `<dl>` `</dl>`. A circled line next to characters is marked by `<cl>` `</cl>`. A wavy line next to characters is marked by `<wl>` `</wl>`.

**Please note:** In old texts, the lines are to the right of characters. In modern texts, the lines may also be to the left of characters. The position to the left or right is not encoded.

##### *Example*

*underlinings to the left and right*

`<sl>`康熙`</sl>``<wl>`寶應縣志`</wl>`人物志： 柏叢桂  
有司寢不行。`<cl>``<sl>`叢桂`</sl>`更相度地形畫圖`</cl>`。



→ For an example with underlinings to the right of characters see section 2.2.2. This example includes `<dl>` for double lines.

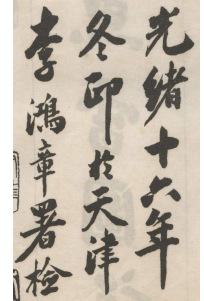
### 3.3.3 Individualized Character Style

A paragraph in an individualized character style is marked by `ics` in the `<p>` tag, i.e. `<p ics>`.

If a character variant in a `<p ics>` paragraph is not in Unicode, resolve it silently, i.e. do not use `< v>` or `< R>`. However, use `<001>` if a character is unknown.

#### Example

`<p ics>`光緒十六年  
冬印於天津  
李鴻章署檢`</p>`



### 3.3.4 Handwritten Text

If a character has been crossed out, mark it by `{ / }`, e.g. `{四/}`. If a character has been inserted, mark it by `{ }`, e.g. `{五}`. If the inserted character replaces a crossed-out character, mark this by `{四/五}`, or `{@/五}` if the crossed-out character is no longer readable.

If there is a line between two consecutive characters indicating that the order should be reversed, mark this by `{ ~ }`, but type the characters in the order as they appear in the text.

If a character variant in handwritten text is not in Unicode, resolve it silently, i.e. do not use `< v>` or `< R>`. However, use `<001>` if a character is unknown.

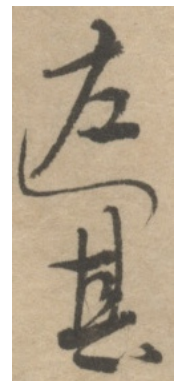
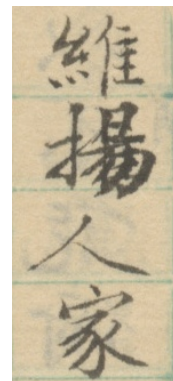
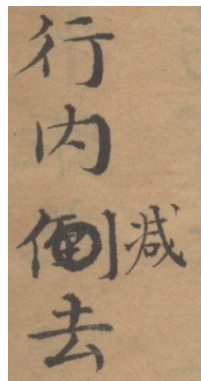
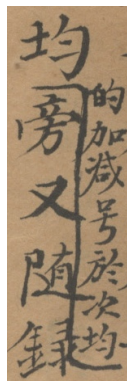
#### Examples

均`{ 的加減號於次均 }`旁又隨錄

行內`{ 倒/減 }`去

維`{ @/揚 }`人家

`{ 左~其 }`



## A Old Radicals That Need to be Marked

**Please note:** This list is only a guideline and may not be complete.

An entry such as 𠂇 𠂇 (𠂇) means that if a character has 𠂇 or 𠂇 as the shape of its radical, and the character with this variant radical is not included in Unicode 5.1, type the standard character and add < R>. Example: <蒼R>.

Example of a radical that does not have to be marked: Do not mark the old shape 𠂇 of the radical 𠂇 because of point a) in the list in section 3.2.2.

*List: old radicals that need to be marked*

𠂇 (𠂇)	𠂇 (𠂇)	𠂇 (𠂇)	𠂇 (𠂇)
臣 (臣)	𠂇 (西)	𠂇 (角)	𠂇 (豕)
𠂇 (走)	辰辰 (辰)	𠂇 (𠂇)	采 (采)
𠂇 (長)	𠂇 (門)	𠂇 (𠂇)	青 (青)
面 (面)	𠂇 (風)	飛 (飛)	𠂇 (食)
𠂇 (首)	𠂇 (馬)	高 (高)	𠂇 (魚)
𠂇 (𠂇)	𠂇 (鹿)	黃 (黃)	黑 (黑)
齊 (齊)	𠂇 (齒)	𠂇 (龜)	𠂇 𠂇 (𠂇)

## B List of All Tags

section	tag	name	tag may contain
2.1	<pb>	page break	page number, a/b
2.1	<rh> </rh>	running head	
2.2.1	<h> </h>	heading	level number
2.2.1	<ti> </ti>	<i>ti</i> 題	
2.2.2	<p> </p>	paragraph	i, x, ics
2.3.1	<tb> </tb>	table	
2.3.1	#	field separator	
2.3.2	<list> </list>	list	
2.3.3	<toc> </toc>	table of contents	
2.4.1	<fig> </fig>	figure	
2.4.1	<cap> </cap>	figure caption	
2.4.1	<desc> </desc>	figure description	
2.4.1	<var> </var>	figure variables	
2.4.2	<stamp> </stamp>	stamp	
2.5.1	<?>	uncertain text	
2.5.1	@, <gap>	unreadable text	
2.5.2	<001>, etc.	unknown character	
3.2.2	< R>	old radical	
3.2.2	< V>	character variant	
3.2.2	< RV>	< R> and < V>	
3.3.1	<sm> </sm>	small characters	
3.3.1	\\	half-column breaks	
3.3.2	<sl> </sl>	single line	
3.3.2	<dl> </dl>	double line	
3.3.2	<cl> </cl>	circled line	
3.3.2	<wl> </wl>	wavy line	
3.3.4	{ / }	crossed-out character	
3.3.4	{ }	inserted character	
3.3.4	{ / }	character replacement	
3.3.4	{ ~ }	reversed characters	